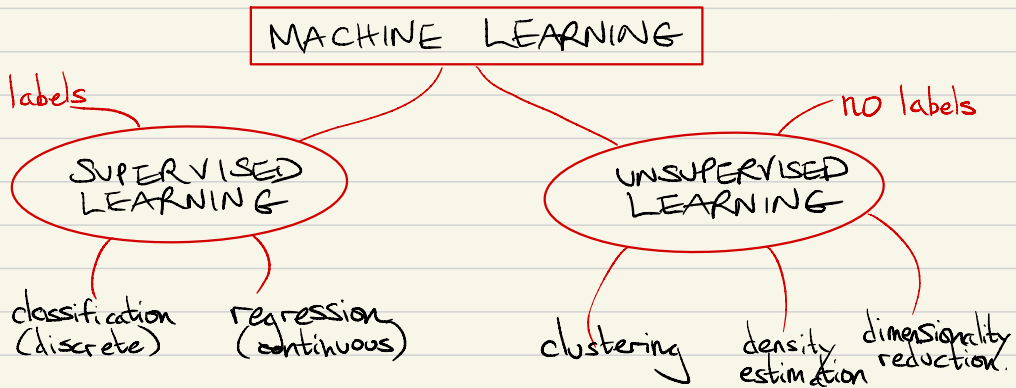


Lecture 13

03/01/2022

Intro to Scikit-Learn



Terminology

* Data recorded in matrix form

$$\hookrightarrow [X] = \underbrace{[N_{\text{samples}}, N_{\text{features}}]}_{\substack{\# \text{ datapoints} \quad \# \text{ attributes}}}$$

* Some data will have labels (sklearn calls these "TARGETS") — stars, galaxies etc.

- Scikit - Learn Workflow

- ① Instantiate an **estimator** object
- ② Fit estimator on **data** and **labels**
- ③ **Predict** new labels
- ④ Find ^{or} **model parameters**

➔ Usually partition dataset into **TRAINING** and **TESTING** sets.

- Supervised Learning — **labels**

- * **Classification** algorithms are trained on labeled data, and used to classify new object features.
- * **Regression** (or "fitting") is the continuous form of classification.

- Unsupervised Learning — **no labels**.

- * Use the data to discover its own labels.
- * Clustering (group similar data), density estimation (find the PDF), dimensionality reduction (find important features).